



CC-BY-NC-SA computer wizard industries 2018

love this?

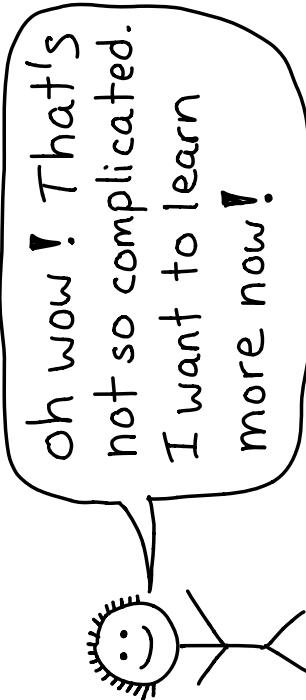
find more awesome zines at

→ [jvns.ca/zines](http://jvns.ca/zines) ←



You're in the right place! This zine has 19 comics explaining important Linux concepts

... 5 minutes later ...



by Julia Evans  
<https://jvns.ca>  
[twitter.com/b0rk](https://twitter.com/b0rk)

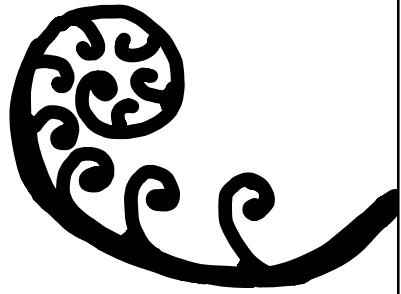
Want to learn more?  
I highly recommend  
this book →

# THE LINUX PROGRAMMING INTERFACE

MICHAEL KERRISK

Every chapter is a readable,  
short (usually 10-20 page)  
explanation of a Linux system.  
I used it as a reference  
constantly when writing  
this zine.

I got it because even though  
it's huge and comprehensive  
(1500 pages!), the chapters  
are short and self-contained  
and it's very easy to pick it  
up and learn something.



# man page sections 22

man pages are split up  
into 8 sections

① ② ③ ④ ⑤ ⑥ ⑦ ⑧

\$ man 2 read

means "get me the man page  
for `read` from section 2"

There's both

→ a program called "read"  
→ and a system call called "read"

so

\$ man 1 read

gives you a different man page from

\$ man 2 read

If you don't specify a section, man will  
look through all the sections & show  
the first one it finds

man page sections

① programs ② system calls

\$ man grep

\$ man sendfile

\$ man ls

\$ man ptrace

③ C functions ④ devices

\$ man printf

\$ man null

\$ man fopen

for /dev/null docs

⑤ file formats ⑥ games

\$ man sudoers

for /etc/sudoers

\$ man proc

files in /proc!

\$ man sl

not super useful.

\$ man s1

is funny if you have

\$ man ls

sl though.

⑦ miscellaneous explains concepts!

\$ man 7 pipe

\$ man apt

\$ man 7 symlink

\$ man chroot

## ♥ Table of contents ♥

unix permissions	4	sockets	10	virtual memory	17
/proc	5	unix domain sockets	11	shared libraries	18
system calls	6	processes	12	copy on write	19
signals	7	threads	13	page faults	20
file descriptors	8	floating point	14	mmap	21
pipes	9	file buffering	15	man page sections	22
memory allocation	16				

# unix permissions

4

There are 3 things you can do to a file

→ **read** → **write** → **execute**

ls -l file.txt shows you permissions  
Here's how to interpret the output:

rw- r-- bork staff  
↑ ↑  
can read & write can read & write  
bork (user) staff (group)

File permissions are 12 bits

setuid setgid  
000 110 110 100  
sticky rw x rwx x  
For files:  
r = can read  
w = can write  
x = can execute  
For directories it's approximately:  
r = can list files  
w = can create files  
x = can cd into & access files

110 in binary is 6

So rw- r-- = 110 100 100  
= 6 4 4

chmod 644 file.txt +  
means change the  
permissions to:  
rw- r-- r--  
Simple!

**setuid** affects executables

\$ ls -l /bin/ping  
rws r-x r-x root root  
this means ping always runs as root  
**setgid** does 3 different unrelated things for executables, directories, and regular files  
it's a long story why???

# mmap

21

What's mmap for?

I want to work with a VERY LARGE FILE but it won't fit in memory

You could try mmap!  
(mmap = "memory map")

load files lazily with mmap

When you mmap a file, it gets mapped into your program's memory  
2TB of virtual memory  
but nothing is ACTUALLY read into RAM until you try to access the memory (how it works: page faults!)

Even if 10 processes

mmap a file, it will only

be read into memory once.

how to mmap in Python

```
import mmap  
f = open("HUGE.txt")  
mm = mmap.mmap(f.fileno(), 0)
```

this won't read the file from disk!  
Finishes ~instantly.

```
print(mm[-1000])
```

this will read only the last 1000 bytes!

anonymous memory maps

→ not from a file (memory set to 0 by default)

→ with MAP\_SHARED, you can use them to share memory with a subprocess!

you too eh? no problem  
I always mmap, so that file is probably loaded into memory already

dynamic linking uses mmap

I need to use libc.so.6  
Standard library  
no problem!

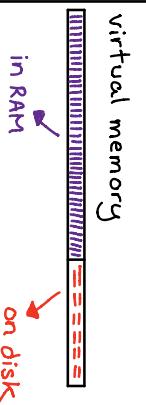
# Page faults

every Linux process has a page table

## \* page table \*

virtual memory address	physical memory address
0x19723000	0x1422000
0x19724000	0x1423000
0x1524000	not in memory
0x1844000	0x1a000 read only

"not in memory" usually means the data is on disk!



Having some virtual memory that is actually on disk is how swap and mmap work

## an amazing directory: `/proc` 5

Every process on Linux has a PID (process ID) like 42.

In `/proc/42`, there's a lot of **VERY USEFUL** information about process 42

### `/proc/PID/fd`

Directory with every file the process has open!

Run `$ ls -l /proc/42/fd` to see the list of files for process 42.

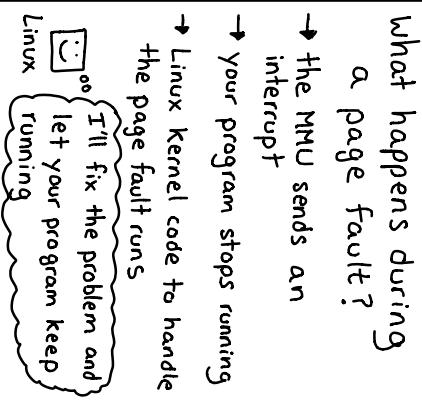
These symlinks are also magic & you can use them to recover deleted files! ♡

some pages are marked as either

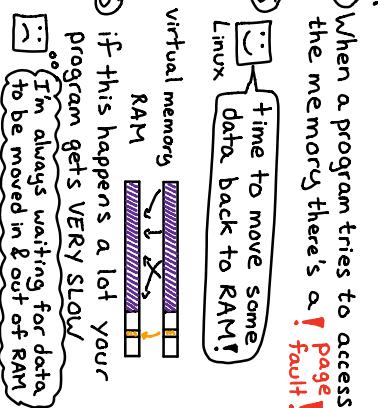
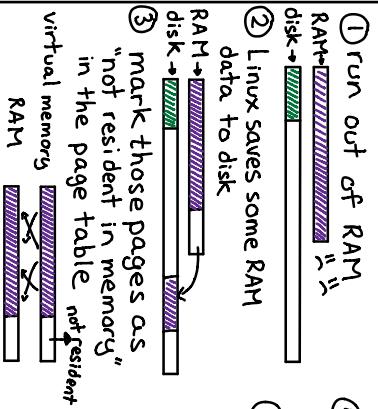
\* read only

\* not resident in memory

when you try to access a page that's marked "not in memory," that triggers a ! page fault!



### how swap works



### `/proc/PID/cmdline`

command line arguments the process was started with

### `/proc/PID/environ`

all of the process's environment variables

### `/proc/PID/exe`

symlink to the process's binary magic: works even if the binary has been deleted!

### `/proc/PID/status`

Is the program running or asleep? How much memory is it using? And much more!

### `/proc/PID/stack`

The kernel's current stack for the process. Useful if it's stuck in a system call

### `/proc/PID/maps`

List of process's memory maps. Shared libraries, heap, anonymous maps, etc.

for more information!

`man proc`

And :more:

Look at

# System calls

<p>The Linux kernel has code to do a lot of things</p> <ul style="list-style-type: none"><li>read from a hard drive</li><li>make network connections</li><li>create new process</li><li>kill a process</li><li>Keyboard drivers</li><li>change file permissions</li></ul>	<p>your program <u>doesn't</u> know how to do those things</p> <ul style="list-style-type: none"><li>"TCP? dude I have no idea how that works"</li><li>"NO I do not know how the ext4 filesystem is implemented I just want to read some files"</li></ul>	<p>and every system call has a number (eg chmod is #90 on x86-64)</p> <p>So what's actually going on when you change a file's permissions is</p> <ul style="list-style-type: none"><li>run syscall #90 program with these arguments</li></ul>
<p>every program uses system calls</p> <ul style="list-style-type: none"><li>I use the 'open' syscall to open files</li><li>me too!</li><li>me three!</li><li>C program</li></ul>		

<p>programs ask Linux to do work for them using <u>=system calls</u>:</p> <ul style="list-style-type: none"><li>please write to this file</li><li>switch to running kernel code</li></ul>	<p>done! I wrote 1097 bytes! &lt;program resumes&gt;</p>	<p>you can see which system calls a program is using with <u>strace</u></p> <ul style="list-style-type: none"><li>\$ strace ls /tmp will show you every system call 'ls' uses!</li><li>it's really fun!</li><li>strace is high overhead</li><li>don't run it on your production database</li></ul>
---	--	--

# copy on write

<p>On Linux, you start new processes using the <u>fork()</u> or <u>clone()</u> system call</p> <p>calling fork gives you a child process that's a copy of you</p> <p>Parent child</p>	<p>copying all that memory every time we fork would be <u>slow</u> and a <u>waste of RAM</u></p> <p>often processes call <u>exec</u> right after <u>fork</u> which means they don't use the parent process's memory basically at all!</p>	<p>when a process tries to write to a shared memory address</p> <ol style="list-style-type: none"><li>① there's a <u>=page fault</u>:</li><li>② Linux makes a copy of the page &amp; updates the page table</li><li>③ the process continues, blissfully ignorant</li></ol> <p>but marks every page as <u>read only</u></p>
---	---	--

<p>programs ask Linux to do work for them using <u>=system calls</u>:</p> <ul style="list-style-type: none"><li>please write to this file</li><li>switch to running kernel code</li></ul>	<p>done! I wrote 1097 bytes! &lt;program resumes&gt;</p>	<p>you can see which system calls a program is using with <u>strace</u></p> <ul style="list-style-type: none"><li>\$ strace ls /tmp will show you every system call 'ls' uses!</li><li>it's really fun!</li><li>strace is high overhead</li><li>don't run it on your production database</li></ul>
---	--	--

# Shared libraries

18

Most programs on Linux use a bunch of C libraries some popular libraries:

`openssl`

(for SSL!)

`sqlite`

(embedded db!)

`libpcre`

(regular expressions!)

`zlib`

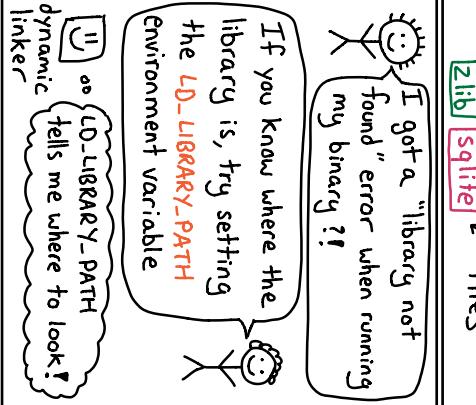
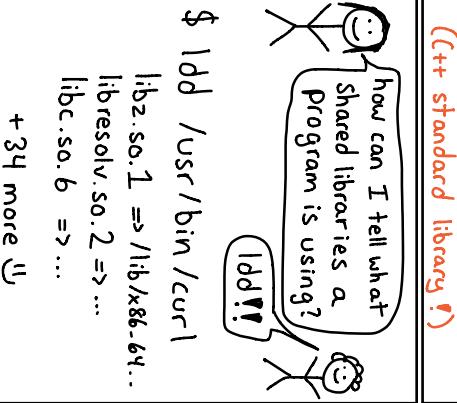
(gzip!)

`libstdc++`

(C++ standard library!)

`sqlite`

(gzipped!)



- ① Link it into your binary big binary with lots of things!  
`[your code] z lib sqlite`
- ② Use separate shared libraries  
`[your code] z lib sqlite` ← all different files
- ③ `LD_LIBRARY_PATH` in executable  
`/etc/ld.so.cache`  
(run `ldconfig -p` to see contents)
- ④ `/lib`, `/usr/lib`

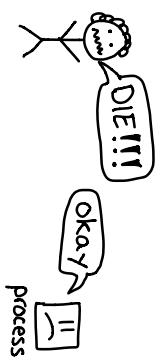
Programs like this:  
`[your code] z lib sqlite` are called "statically linked"

and programs like this:  
`[your code] z lib sqlite` are called "dynamically linked"

# Signals

7

If you've ever used `KILL` you've used signals



the Linux kernel sends your process signals in lots of situations

- `your child terminated`
- `that pipe is closed`
- `illegal instruction`
- `the timer you set expired`
- `Segmentation fault`

You can send signals yourself with the `kill` system call or command

`SIGINT` `ctrl-C` } various levels of  
`SIGTERM` `kill` } "die"  
`SIGKILL` `kill -q`  
`SIGHUP` `kill -HUP`  
often interpreted as "reload config", eg by nginx

Signals can be hard to handle correctly since they can happen at ANY time

SURPRISE! another signal!  
Handling a signal

- |  |   |
|--|---|
| Every signal has a default action, one of: | <code>ignore</code>                               |
|  | <code>KILL</code> process                         |
|  | <code>KILL</code> process AND make core dump file |
|  | <code>STOP</code> process                         |
|  | <code>RESUME</code> process                       |

- |                                   |                              |
|-----------------------------------|------------------------------|
| SIGTERM (terminate) up then exit! | ok! I'll clean up then exit! |
| SIGSTOP & SIGKILL                 | can't be ignored             |
| got →                             | process                      |

- |                        |    |
|------------------------|----|
| <code>SIGKILLED</code> | !! |
|------------------------|----|

## file descriptors

8

file descriptors can refer to:

- files on disk
- pipes
- sockets (network connections)
- terminals (like xterm)
- devices (your speaker! /dev/null!)
- LOTS MORE (eventfd, inotify, signalfd, epoll, etc etc)

**not EVERYTHING on Unix is a file, but lots of things are**

**lsof** (list open files) will show you a process's open files

```
$ lsof -p 4242
  PID  NAME
  0  /dev/pts/pty1
  1  /dev/pts/pty1
  2  pipe:29174
  3  /home/bark/awesome.txt
  5  /tmp/
```

FD is for file descriptor

Unix systems use integers to track open files

Open foo.txt

Okay! that's file #7 for you.

these integers are called **file descriptors**

When you read or write to a file / pipe / network connection you do that using a file descriptor

connect to google.com

ok! fd is 5!

OS write GET / HTTP/1.1 to fd #5 done!

Behind the scenes:

```
f = open("file.txt")
f.readlines()
```

stdin → 0      stdout → 1      stderr → 2

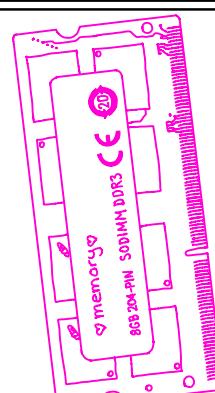
"read from stdin"      "means"      "read from the file descriptor 0"

"read from the file descriptor 0" → could be a pipe or file or terminal!

# virtual memory

7

your computer has physical memory



Physical memory has addresses 0 - 8GB but when your program references an address like 0x5c69a2a2 ↑ that's not a physical memory address! It's a **virtual** address

every program accesses a virtual address

I'm accessing 0x21000 CPU

I'll look that up in the **page table** and then access the right physical address

a "page" is a 4kB or sometimes bigger chunk of memory

PID	virtual addr	physical addr
1971	0x20000	0x192000
2310	0x20000	0x228000
2310	0x21000	0x9788000

every program has its own virtual address space

0x129520 → "puppies" program 1

0x129520 → "bananas" program 2

here's the address of process 2950's page table Linux

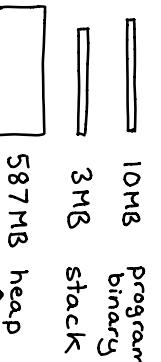
MMU hardware

thanks I'll use that now!

# memory allocation

16

your program has memory



the heap is what your allocator manages

your memory allocator's interface

`malloc(size_t size)` allocate `size` bytes of memory & return a pointer to it  
`free(void* pointer)` mark the memory as unused (and maybe give back to the OS)  
`realloc(void* pointer, size_t size)` ask for more /less memory for `pointer`  
`calloc(size_t members, size_t size)` allocate array + initialize to 0

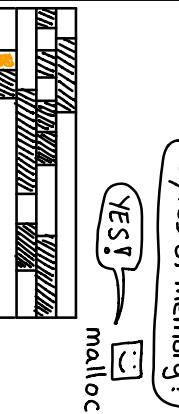
Your memory allocator (`malloc`) is responsible for 2 things.

**THING 1:** keep track of what memory is used / free



`malloc` tries to fill in unused space when you ask for memory

`malloc` can I have 512 bytes of memory?



your new memory

oh no I'm being asked for 40MB and I don't have it

can I have 60MB more?  
here you go!

OS

`malloc` isn't :magic: it's just a function!

You can always:  
→ use a different `malloc` library like `jemalloc` or `tcmalloc` (easy!) → implement your own `malloc` (harder)

## pipes

q

Sometimes you want to send the output of one process to the input of another

\$ ls | wc -l

53  
53 files!

a pipe is a pair of 2 magical file descriptors  
① pipe input `IN` → `our` pipe  
② pipe output  
ls → wc → stdout

When `ls` does `write(IN, "hi")` `WC` can read it!  
read(`our`) → "hi"

Pipes are one-way. → You can't write to `our`.

Linux creates a buffer for each pipe

ls → buffer → wc  
IN → data waiting to be read our

If data gets written to the pipe faster than it's read, the buffer will fill up. `IN` → `out`

When the buffer is full, writes to `IN` will block (wait) until the reader reads. This is normal & ok!)

what if your target process dies?

ls → wc  
\*\* → \*\*

If `wc` dies, the pipe will close and `ls` will be sent `SIGPIPE`. By default `SIGPIPE` terminates your process.

### named pipes

\$ mkfifo my-pipe  
This lets 2 unrelated processes communicate through a pipe!

```
f=open("./my-pipe")
f.write("hi!\n")
f.readline() < "hi!"
```

# Sockets

10

Networking protocols are complicated

Unix systems have an API called the "socket API" that makes it easier to make network connections (Windows too! 😊)

you don't need to know how TCP works, Unix I'll take care of it!

what if I just want to download a cat picture

TCP/IP Illustrated Volume 1 Stevens 600 pages

here's what getting a cat picture with the socket API looks like:

- ① Create a socket
- ② Connect to an IP/port  
Connect(fd, 12.15.14.15:80)
- ③ Make a request  
write(fd, "GET /cat.png HTTP/1.1")
- ④ Read the response  
cat-picture = read(fd ...)

AF-INET?

What's that?

AF-INET means basically "internet socket": it lets you connect to other computers on the internet using their IP address.

The main alternative is AF-UNIX ("unix domain socket") for connecting to programs on the same computer

Every HTTP library uses sockets under the hood

\$ curl awesome.com  
\$ curl awesome.com  
Python: requests.get("yay.us")  
oh, cool, I could write a HTTP library too if I wanted. Neat!  
\* SO MANY edge cases though! 😊

## file buffering

15

I/O libraries don't always call `write` when you print

```
printf("I ❤️ cats");
```

☞ I'll wait for a newline printf before actually writing

This is called **buffering** and it helps save on syscalls

On Linux you write to files & terminals with a system call called **write**

```
please write "I ❤️ cats" to file #1 (stdout)
```

☞ Okay! Linux

? ? ? I printed some text but it didn't appear on the screen. why?? time to learn about **flushing**!

To force your IO library to write everything it has in its buffer right now, call **flush**!

☞ I'll call `write` right away!!

3 kinds of buffering (defaults vary by library)

- ① None. This is the default for `stderr`
- ② Line buffering. (write after newline). The default for `terminals`.
- ③ "full" buffering. (write in big chunks). The default for `files` and `pipes`.

☞ no seriously, actually write to that pipe please

# floating point

14

a double is 64 bits.

sign exponent fraction  
 $\downarrow$   
 101011 101011 101011 101011  
 101011 101011 101011 101011

$$\pm 2 \times 1.\text{frac}$$

That means there are  $2^{64}$  doubles.  
 The biggest one is about  $2^{1023}$

doubles get farther apart as they get bigger  
 between  $2^n$  and  $2^{n+1}$  there are always  $2^{52}$  doubles,  
 evenly spaced  
 that means the next double after  $2^{60}$  is  $2^{60} + 64 = \frac{2^{60}}{2^{52}}$

$$2^{52} + 0.2 = 2^{52}$$

$\leftarrow$  (the next number after  $2^{52}$  is  $2^{52} + 1$ )

$$1 + \frac{1}{2^{54}} = 1$$

$\leftarrow$  (the next number after 1 is  $1 + \frac{1}{2^{52}}$ )

$2^{2000} = \text{infinity}$

$\leftarrow$  infinity is a double

infinity - infinity = nan

Javascript only has doubles (no integers!)

$$> 2^{**} 53 \\ 9007199254740992 \\ > 2^{**} 53 + 1 \\ 9007199254740992$$

$\uparrow$  same number! uh oh!

doubles are scary and their arithmetic is weird

they're very logical! just understand how they work and don't use integers over  $2^{53}$  in Javascript

## Unix domain sockets !!

unix domain sockets are files.

```
$ file mysock.sock
socket
```

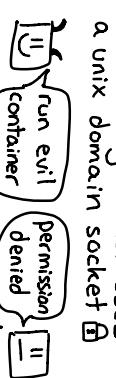
the file's permissions determine who can send data to the socket

### advantage 1

lets you use file permissions to restrict access to HTTP/ database services!

chmod 600 secret.sock

This is why Docker uses a unix domain socket



Linux

### advantage 2

UDP sockets aren't always reliable (even on the same computer).  
 unix domain datagram sockets are reliable!  
 And won't reorder!

I can send data and I know it'll arrive

they let 2 programs on the same computer communicate.  
 Docker uses Unix domain sockets, for example!

Here you go!

### advantage 3

There are 2 kinds of unix domain sockets:  
 stream like TCP! Lets you send a continuous stream of bytes  
 datagram like UDP! Send discrete chunks of data

here's a file I downloaded from sketchy.com

video decoder  
 Sandboxed process

# what's in a process?

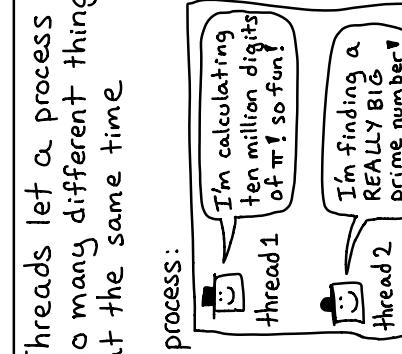
PID	USER and GROUP process #129 reporting for duty!	SIGNAL HANDLERS like PATH! you can set them with: \$ env A=val ./program
WORKING DIRECTORY	Relative paths (./blah) are relative to the working directory! <code>chdir</code> changes it.	PARENT PID PID 1 (init) ↓ PID 147 ↑ PID 129
MEMORY	heap! stack! shared libraries! the program's binary! mmaped files!	CAPABILITIES I have CAP_PTRACE well I have CAP_SYS_ADMIN
PID	who are you running as? julia!	ENVIRONMENT VARIABLES like PATH! you can set them with: \$ env A=val ./program
OPEN FILES	COMMAND LINE ARGUMENTS see them in /proc/PID/cmdline	NAMESPACES I'm in the host network namespace I have my own namespace! container process

12

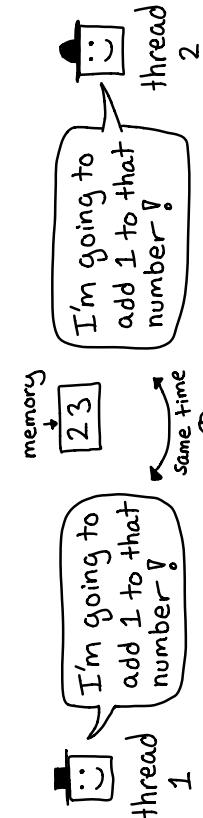
13

# threads

Threads let a process do many different things at the same time	and they share code
process:	calculate-pi: find-big-prime-number
thread 1 I'll write some digits of π to 0x129420 in memory	but each thread has its own stack and they can be run by different CPUs at the same time
thread 2 oh that's where I was putting my prime numbers	CPU 1 CPU 2 π thread primes thread



Sharing memory can cause problems (race conditions!)



RESULT: 24 ↪ WRONG.  
Should be 25!

sharing data between threads is very easy. But it's also easier to make mistakes with threads

You weren't supposed to CHANGE that data!  
thread 1

why use threads instead of starting a new process?  
→ a thread takes less time to create